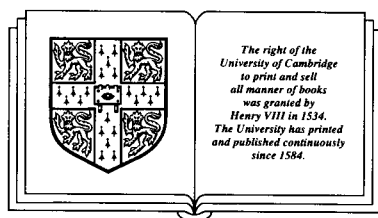

FORMAL SEMANTICS AND PRAGMATICS FOR NATURAL LANGUAGE QUERYING

JAMES CLIFFORD

New York University



CAMBRIDGE UNIVERSITY PRESS

Cambridge

New York Port Chester Melbourne Sydney

Published by the Press Syndicate of the University of Cambridge
The Pitt Building, Trumpington Street, Cambridge CB2 1RP
40 West 20th Street, New York, NY 10011, USA
10 Stamford Road, Oakleigh, Melbourne 3166, Australia

© Cambridge University Press 1990

First published 1990

Printed in Great Britain at the University Press, Cambridge

Library of Congress cataloguing in publication data available

British Library cataloguing in publication data available

ISBN 0 521 35433 1

Contents

1	Introduction	1
2	Montague Semantics	5
2.1	Overview of the Approach	5
2.2	The Language IL_s	6
2.3	Semantics of IL_s	7
2.4	IL_s and Montague's IL	9
2.5	Informal Discussion of IL_s	11
3	The HRDM Model	17
3.1	Motivation for Historical Databases	17
3.2	Informal Presentation of HRDM	20
3.3	Lifespans	23
3.4	Intuitive Presentation of Historical Databases	28
3.5	Historical Relations in HRDM	37
4	Intensional Logic and HRDM Model	45
4.1	Introduction	45
4.1.1	Introduction	45
4.1.2	The IL_s Language Defined by an HRDB Scheme	48
4.2	The Intensional Model Induced by a Database Instance	54
4.2.1	The Interpretation of the Non-Logical Constants	59
4.3	Informal Discussion of IL_s and HRDB	63
4.3.1	Domains and values	63
4.3.2	Attributes	63
4.3.3	Tuples	64
4.3.4	Data Dependencies and Constraints	65
4.3.5	Queries	67
5	Overview of English Query Language QE-III	69
5.1	Introduction	69
5.2	Preliminaries	71
5.2.1	Individual Concepts vs. Entities	71
5.2.2	Verbs	71

5.3	The Problems of Tense and Time	72
5.3.1	Intervals or States?	72
5.3.2	Sentential vs. Verb-phrasal Temporal Operators	76
5.4	Questions	78
5.4.1	Introduction	78
5.4.2	Database Questions	80
5.5	The QE-III Theory of Questions	84
5.5.1	Introduction	84
5.5.2	Yes-No Questions	84
5.5.3	WH-Questions	85
5.5.4	Temporal Questions	95
5.5.5	Pragmatics or Semantics?	96
5.5.6	Conclusions	100
5.6	Related Work	101
5.6.1	Introduction	101
5.6.2	Karttunen	101
5.6.3	Bennett and Belnap	103
5.6.4	Hausser and Zaefferer	105
5.6.5	Scha and Gunji	106
6	Formal Definition of QE-III	109
6.1	Introduction	109
6.2	The Syntax of QE-III	110
6.2.1	The Categories	110
6.2.2	The Syntactic Rules of Formation	112
6.3	The Semantics of QE-III	120
6.4	The Pragmatics of QE-III	124
6.5	A QE-III Fragment Schema	127
7	Examples from the QE-III Fragment	133
7.1	Introduction	133
7.2	PTQ-like Examples from the QE-III Fragment	134
7.3	Temporal Reference in QE-III	137
7.4	Questions in QE-III	146
7.5	Miscellaneous Features of QE-III	160
7.6	Query Evaluation in Practice	166
7.7	Conclusion	170
8	Summary and Conclusions	173
8.1	Summary	173
8.1.1	The Historical Relational Database Model	173
8.1.2	The Language QE-III	174
8.2	Future Work	175
8.2.1	Time and Databases	175

8.2.2	Formalized Natural Query Languages	177
-------	--	-----

List of Figures

3.1	Relation emprel on Scheme EMPREL	20
3.2	Three Intensions	23
3.3	A Relational Database Instance	24
3.4	The Three Dimensions of a Historical Database	25
3.5	One Lifespan Associated with Entire Database	25
3.6	Lifespans Associated with Each Relation	26
3.7	Lifespans Associated with Each Relation Tuple	26
3.8	Relational Database Schema	27
3.9	Lifespan of Attribute DAILY-TRADING-VOLUME	27
3.10	Interaction of Tuple and Attribute Lifespans	28
3.11	Tuple and Attribute Lifespans	29
3.12	Three Static Relation Instances	30
3.13	Three Extended Static Relations	31
3.14	Historical Relations as Three-Dimensional	32
3.15	<i>Completed</i> Relations	34
3.16	A Historical Relation as a 3-D Structure	35
3.17	Completed Relation	36
3.18	Example Historical Relational Database in HRDM	40
3.19	Levels in Historical Relational Data Model	41
4.1	Relation emprel	49
4.2	Relation deptrel	49
4.3	Relation itemrel	50
4.4	Relation salesrel	50
4.5	Attribute Values Specified at Points in Time	56
4.6	Attribute Value as a Step Function	57
5.1	Derivation of “John worked yesterday.”	75
5.2	Translation of “John worked yesterday.”	75
5.3	Time Line Consistent with Logical Translation	76
5.4	PTQ Derivation of “John walks.”	76
5.5	Dowty’s Derivation of “John walked.”	77
5.6	QE-III Derivation of “John worked.”	78
5.7	One Possible Derivation of “Who manages every employee?”	82

5.8	A Possible Translation of "Who manages every employee?"	82
5.9	Another Possible Derivation of "Who manages every employee?" . . .	82
5.10	A Possible Translation of "Who manages every employee?"	83
5.11	Similarity of Question Terms and Terms in Subject Position	86
5.12	WH-Q Movement in Object Position	86
5.13	Similarity of Semantics of Terms and WH-Terms	87
5.14	Translation of "Peter manages the shoe department."	87
5.15	Translation of "Who manages the shoe department."	88
5.16	Two PTQ Derivations of "John walks."	89
5.17	Pronominal Co-reference	89
5.18	Co-reference in Questions	90
5.19	Blocked Analysis	91
5.20	Translation of Lower Constituents in Blocked Analysis	91
5.21	Multiple "who" Semantics	92
5.22	A Multiple "who" Semantics Treatment of "Who manages what?" . .	92
5.23	An Infinite WH-Rule Semantics Treatment of "Who manages what?"	94
5.24	Simultaneous Substitution of WH-Terms	94
5.25	Wide Scope of "when"	96
5.26	Translation of "When did John work?"	97
5.27	Variables in Bennett/Belnap Theory	104
5.28	Questions in QE-III and H-Z	106
6.1	Variables Used in QE-III Translations	120
7.1	Variables Used in QE-III Translations	134
7.2	QE-III Derivation of "John manages Mary."	135
7.3	QE-III Translation of "John manages Mary."	136
7.4	QE-III Derivation of "Peter earned 25K in 1978."	137
7.5	QE-III Translation of "Peter earned 25K in 1978."	138
7.6	An Incorrect Derivation of "Peter earned 25K in 1978."	139
7.7	Translation Corresponding to Incorrect Derivation of "Peter earned 25K in 1978."	139
7.8	Time Line Consistent with Incorrect Derivation	140
7.9	Another Incorrect Derivation of "Peter earned 25K in 1978."	140
7.10	Translation Corresponding to Another Incorrect Derivation of "Peter earned 25K in 1978."	141
7.11	Time Line Consistent with Second Incorrect Derivation	141
7.12	QE-III Derivation of "Peter manages an employee such that he earned 30K."	142
7.13	QE-III Translation of "Peter manages an employee such that he earned 30K."	143
7.14	QE-III Derivation of "John worked before Mary worked."	143
7.15	QE-III Translation of "John worked before Mary worked."	144
7.16	QE-III Derivation of "Rachel worked before yesterday."	145

7.17 QE-III Translation of "Rachel worked before yesterday."	145
7.18 QE-III Derivation of "Who managed Rachel?"	146
7.19 QE-III Translation of "Who managed Rachel?"	147
7.20 Incorrect Derivation of "Who managed Rachel?"	148
7.21 Translation Corresponding to Incorrect Derivation of "Who managed Rachel?"	148
7.22 QE-III Derivation of "Who manages which employees?"	149
7.23 QE-III Translation of "Who manages which employees?"	150
7.24 QE-III Derivation of "What does who supply to whom?"	151
7.25 QE-III Translation of "What does who supply to whom?"	152
7.26 QE-III Derivation of "Who works for a department such that it sells shoes?"	153
7.27 QE-III Translation of "Who works for a department such that it sells shoes?"	154
7.28 QE-III Derivation of "Is it the case that Peter earns 30K?"	155
7.29 QE-III Derivation of "Does Peter earn 30K?"	155
7.30 QE-III Translation of "Is it the case that Peter earns 30K?"	156
7.31 QE-III Translation of "Does Peter earn 30K?"	156
7.32 QE-III Derivation of "When did Peter earn 25K?"	157
7.33 QE-III Translation of "When did Peter earn 25K?"	157
7.34 QE-III Derivation of "When did who manage whom?"	157
7.35 QE-III Translation of "When did who manage whom?"	158
7.36 QE-III Derivation of "When and to whom did company A supply item B yesterday?"	158
7.37 QE-III Translation of "When and to whom did company A supply item B yesterday?"	159
7.38 QE-III Derivation of "Who is Peter's manager?"	160
7.39 QE-III Translation of "Who is Peter's manager?"	161
7.40 QE-III Derivation of "Who is a manager of Peter?"	162
7.41 QE-III Translation of "Who is a manager of Peter?"	163
7.42 QE-III Derivation of "Who has Peter as manager?"	164
7.43 QE-III Translation of "Who has Peter as manager?"	164
7.44 QE-III Derivation of "Who sells item 37?"	165
7.45 QE-III translation of "Who sells item 37?"	165

INTRODUCTION

It is difficult to imagine a successful semantic theory that does not include time as an integral component. Yet all of today's major data models – models which purport to provide a general theory on how to represent information for convenient and rapid storage and retrieval on digital computers – completely ignore this essential aspect of semantics. In this work we examine the connection between two areas of semantics, namely the semantics of databases and the semantics of natural language, and link them together via a common view of the semantics of time. In the first part we argue that an essential ingredient for the success of efforts to incorporate more *real world* semantics into database models is a coherent theory of the semantics of time. We describe such a database theory, and then proceed to present a formally defined English database query language whose semantic theory makes explicit reference to the notion of denotation with respect to a moment of time.

The idea that time might be an important consideration in providing an enriched database semantics is not new to this work, but it is nonetheless a relatively recent concern of database research. [Bub77], [Ser80], [Klo81], [And81], [AM82], [CW83], [Sno84], [GV85], [CC87], [GY88], and [CC88a] are among the many works that have lately investigated ways in which time might be added to a database model. Two recent surveys of the literature on time and databases, [BADW82] and [Sno86], and a recent international conference, [TAI87], provide excellent references to the growing interest and literature in this field. Our own approach (first presented in [Cli82a] and [CW83] and further investigated in [CC87]) owes a debt to the philosophical tradition represented in such works as [Car47], [Pri67], and [RU71], and particularly in the works of the logician Richard Montague [Mon74]. This tradition makes a case for a far more pervasive theory of the importance of time in a theory of meaning than one which asserts that time is just one among many equally important aspects.

In Chapter 2 we present a brief introduction to the Montague Semantics (MS) framework, including a definition of the intensional logic IL_s . Then in Chapter 3 we present a formal definition of the Historical Relational Database Model (HRDM), our formal incorporation of a semantics for time within the context of the relational

database model ([Cod70]). Unlike much of the recent database work involving time, we do not attempt to incorporate any notions of how time is encoded (e.g., in the Julian calendric system.) Such theories are important, but it seems that before exploring some of the minutiae of time a theory that captures more of its essence is needed. Theories concerned with months of the year or days of the week are best discussed at the level of interfaces or data representation schemes, and not at the level of a basic semantic theory. We therefore present a very general theory, one which ascribes only the simplest and most intuitive properties to time (order and density) and defines a very simple relationship between time and the other elements of the relational database model.

Our emphasis, as discussed in Chapters 3 and 4, is on the centrality of time with respect to database modelling, and we argue that time is central to our understanding of database semantics. We present the point of view that it is not simply expressions such as “previous employee” or “salary increase” that require reference to the notion of time for their understanding, but rather that every aspect of database theory is understood better when its relationship to the phenomenon of time is considered. In particular, a new insight into the nature of the database distinction between *key* and *non-key* attributes is provided by an understanding of their relationship to time.

In order to make a concrete presentation of these ideas, we present them within the context of two well-known database models, the relational model [Cod70] and the entity-relationship model [Che76]. We adopt the relational model both because it is a well-studied and formalized database model, and because it is increasingly being used as a model for implemented systems. The entity-relationship model is used not simply because of its growing popularity as a tool for modelling database semantics, but also because its ontological theory – that the world consists of entities and relationships among them – is in close accord with the ontology of the philosophical logic tradition.

In the latter part of this work we argue that a fully formalized account can successfully be given for both the syntax and the interpretation of an English database query language. Moreover this account is independent of any performance model of how a processor (for example, a computer) might go about understanding such a language. The language theory that we present, with syntax paired with semantics, is once again a direct outgrowth of the work of Richard Montague, who argued quite seriously in a number of papers for the contention that there is “no important theoretical difference ... between formal and natural languages” ([Mon70a, p.188 in [Mon74]].) Working within the framework of Montague’s syntactic and semantic theory we present a formalized fragment of English questions designed for the purposes of database querying.

Numerous systems for providing natural language access to databases have been described in the literature, including [WKN81], [Wal78], [Har78], and [HSSS78]. While these systems are dissimilar in a number of different respects, they all share what to us is the same defect, namely the lack of any fundamental formal theory of the semantics of the database or the semantics of the English query language.

We view the development of these and other such systems as belonging to the

first phase in the development of a formal theory of the semantics of database and of database querying, much as the early years in the design of computer languages such as FORTRAN were part of the first phase in the development of a theory of programming-language semantics. It awaited the impact of formal language theory, coupled with a theory of syntax-directed translation, for the area of programming language theory to be born. To complete an analogous development in the area of natural language querying would require the impact of formal language theory and a theory that coupled the syntax and the semantics of English. Many linguists today believe that Montague's theory of universal grammar [Mon70c] is the first successful attempt at formalizing such a uniform syntactic and semantic theory of natural language. We believe that some such formal theory of a query language is an important first step towards the development of provably correct and reliable natural language processing systems. For inherent in the notion of program *correctness* is the concept of a standard against which a program is to be judged.

In Chapter 5 we present an informal overview of a fragment of English for database querying which we call QE-III. We discuss the kinds of properties and abilities that a database query language in English should possess. Principal among these are an account of question semantics that possesses close analogs in database theory, an account of the semantics of multiple-WH questions, an account of the semantics of time, and a grammar that is conducive to a computer implementation. After examining a number of partial solutions to these problems, we introduce the notion of pragmatics as an additional formal component of our language's theory. We argue that assigning to the pragmatic component the task of providing a representation for the answer(s) to a question is both appropriate and elegant. Finally we discuss several other recent attempts at developing a formal theory of questions.

After this informal presentation we provide in Chapter 6 a formal definition of the query fragment QE-III as a Montague Grammar. This fragment represents a simplification of the semantic theory of Montague's fragment, presented in [Mon73] and known in the literature as PTQ, that offers a natural correspondence to the semantics of queries in a database context. An excellent introduction to the formal semantic approach to linguistic analysis which has come to be known as Montague Semantics can be found in [DWP81]. The fragment is provided with a formal syntax, semantics, and pragmatics, each component designed with the database application in mind. These components of the QE-III language are based upon the language IL_s introduced in Chapter 2. The inclusion of a formal pragmatic component is an interesting extension to the traditional conception of a Montague Grammar. Among the major extensions to the PTQ fragment embodied in QE-III are the inclusion of time-denoting expressions and temporal operators, an analysis of verb meanings into primitive meaning units derived from the database schema, and of course the inclusion of certain forms of direct questions. These extensions, and the semantics with which they are provided, are motivated by the ultimate goal of database access, but they are equally interesting in their own right. The syntactic theory presented is in some cases admittedly naive, for we have been primarily interested in getting the

interpretation right.

A more complete discussion of the details of the features of the QE-III fragment is presented in Chapter 7, where numerous example derivations and translations of typical database queries are examined and discussed. We also discuss briefly how QE-III can be adapted to different database domains, and how the logical translations of the English queries expressed in QE-III can be translated into a data manipulation language like the historical relational algebra of [CC87]. The chapter concludes with a discussion of some of the limitations of the fragment and of some possibilities for further extensions. Finally, we conclude in Chapter 8 with a review of the major ideas of the book and with a discussion of future work that this research suggests.